

A remark on equivalent Rosser sentences

Christopher von Bülow*

University of Constance, Germany

Received 13 October 2006; received in revised form 8 October 2007; accepted 8 October 2007

Available online 26 November 2007

Communicated by S.N. Artemov

Abstract

An oversight in Guaspari and Solovay's "Rosser sentences" [D. Guaspari, R.M. Solovay, Rosser sentences, *Annals of Mathematical Logic* 16 (1) (1979) 81–99] is pointed out and emended. It concerns the premisses of their proof that there are standard proof predicates all of whose Rosser sentences are provably equivalent. The result holds up, but the premisses mentioned in the paper have to be strengthened somewhat.

© 2007 Elsevier B.V. All rights reserved.

Keywords: Proof predicate; Rosser sentence; Provability logic

1. Introduction

In 1979 David Guaspari and Robert Solovay published a joint paper titled "Rosser sentences" [1].¹ In this paper, they introduce formal systems of propositional modal logic which capture the provability logic of Rosser sentences, i.e., of sentences in the language \mathcal{L}_{PA} of Peano arithmetic which assert: "there is a disproof of me that occurs before any proof of me."

For these formal systems semantical and arithmetical completeness theorems hold, which are formulated and proved by Guaspari and Solovay in Sections 3–5 of their paper. They then go on to apply their results in Parts A and B of Section 6. Thus, in Subsection 6.A they show that there are standard proof predicates² *not* all of whose Rosser sentences are provably equivalent. Finally, in Subsection 6.C, they prove that there also are standard proof predicates (SPP's) all of whose Rosser sentences *are* provably equivalent.

To this purpose, they define an exotic proof predicate which utilizes a list of its own Rosser sentences, the 'Rosser list', in deciding how to behave. For this reason I call the new proof predicate "LTh" — short for "(Rosser) list theorem". The formula $\text{LTh}(f)$ will be of the form $\exists 1 \text{LPF}(1, f)$, where " $\text{LPF}(1, f)$ " stands for " 1 is a 'list proof'

* Corresponding address: Universität Konstanz, Philosophische Fakultät, Fach D 21, D-78457 Konstanz, Germany.

E-mail address: Christopher.von.Buelow@uni-konstanz.de.

URL: <http://www.uni-konstanz.de/FuF/Philo/Philosophie/philosophie/?/88-0-Christopher-von-Buelow.htm>.

¹ The present paper should be read along with [1] to supply the context not given here.

² The concept of a standard proof predicate is a generalization of Gödel's $\text{Bew}(f)$, which says: "there is a **PA**-proof the last line of which is f "; see Definition 1.

for \mathbf{f} ".³ The formula $\text{LPF}(1, \mathbf{f})$ is a pterm (pseudo-term) with respect to \mathbf{f} [2, p. 24], i.e., \mathbf{PA} proves that for every 1, there is one and only one \mathbf{f} such that $\text{LPF}(1, \mathbf{f})$. (The pterm $\text{LPF}(1, \mathbf{f})$ corresponds to the recursive function f in [1, pp. 97–98].)

In defining $\text{LPF}(1, \mathbf{f})$, Guaspari and Solovay start with an arbitrary SPP $\text{Th}(\mathbf{f})$ having a certain additional property (+), given in [1, p. 97]. For example, $\text{Th}(\mathbf{f})$ might be the usual proof predicate $\text{Bew}(\mathbf{f})$. From $\text{Th}(\mathbf{f})$ we obtain a formula $\text{ThNum}(m, \mathbf{f})$ which describes a numeration of the \mathbf{PA} -theorems w.r.t. $\text{Th}(\mathbf{f})$. (In the case of $\text{Bew}(\mathbf{f})$ this could be a numeration replicating the sequence given by the size of the Gödel numbers of corresponding proofs.) Since $\text{ThNum}(m, \mathbf{f})$ is a pterm, we can adopt a convention from [2, p. 24] and write " $\text{thNum}(m)$ " to stand for the \mathbf{f} with $\text{ThNum}(m, \mathbf{f})$, as if " thNum " were a new one-place function letter. (This pterm corresponds to the recursive function g in [1, pp. 97–98].)

The new proof predicate $\text{LTh}(\mathbf{f})$ works according to the following algorithm⁴: In the beginning we let LTh output whatever thNum outputs; that is, we let $\text{LPF}(0, \text{thNum}(0))$, $\text{LPF}(1, \text{thNum}(1))$, and so forth. While this is going on we also check each formula that is output whether it is by any chance one which says, for some sentence ρ :

ρ is a Rosser sentence with respect to $\text{LTh}(\mathbf{f})$,

i.e., whether it is of the form,

$$\rho \leftrightarrow \exists y \left[\text{LPF}(y, \ulcorner \neg \rho \urcorner) \wedge (\forall z \leq y) \neg \text{LPF}(z, \ulcorner \rho \urcorner) \right].$$

If so, we add ρ to the Rosser list — or at least we do so ordinarily. However, we don't want the list to contain sentences ρ_1, ρ_2 for which $\rho_2 = \neg \rho_1$, because then we would have to keep track of which of them comes earlier in the list. In order to avoid this, we *refrain* from adding the Rosser sentence ρ to the Rosser list if the list already includes either $\neg \rho$ or some ρ' such that $\rho = \neg \rho'$.

While LTh outputs Th -theorems and, slowly, the Rosser list is being built up, we also watch whether one of these Th -theorems is (i) a Rosser sentence w.r.t. $\text{LTh}(\mathbf{f})$, or (ii) the negation of one. (Actually, we merely look whether it is one of the sentences on the Rosser list in its current state, or the negation of such a sentence.) When this happens, Guaspari and Solovay say: 'the bell rings.'

In case (i), Th outputs an LTh -Rosser sentence. We let LTh reproduce this behavior by making LTh immediately output the whole Rosser list and, afterwards, all sentences of $\mathcal{L}_{\mathbf{PA}}$ in some arbitrary order.⁵ Assuming that LTh hasn't yet output the *negation* of any of the sentences on the list, all list members are thus output ('proved') by LTh before their negations are. But as Rosser sentences they all assert the converse. Thus they are all false, and therefore equivalent.

In case (ii), i.e., if Th outputs the negation of a Rosser sentence, we make LTh immediately output the negations of all the sentences on the list and then (as before) all sentences of $\mathcal{L}_{\mathbf{PA}}$. Assuming that LTh hasn't yet output any of the Rosser-listed sentences themselves, their negations are thus being 'proved' before they themselves are. This is just what they say, and so they are all true and hence equivalent.

This sketch is of course far from sufficient as a *proof* that all Rosser sentences w.r.t. $\text{LTh}(\mathbf{f})$ are provably equivalent; it is merely intended to provide some conception of the functioning of the new proof predicate $\text{LTh}(\mathbf{f})$, and to hint at the direction which a proof would have to take. Here, it is enough to say that a proof is given in [1, pp. 97–98], but that the condition (+) offered in the paper is slightly too weak to support the proof presented for the central Lemma 6.3 [1, p. 98].⁶ The point of this paper is to emend the condition and to show — in a little more detail than in the original paper — how Lemma 6.3 is to be proved.

2. How to prove Lemma 6.3

Here is Guaspari and Solovay's definition of a standard proof predicate [1, p. 83]:

³ I employ typewriter font to set off object-language variables.

⁴ Cf. [3, [24.1], [24.2]]. The proof predicate is obtained via the diagonal lemma from a formalization of the following description.

⁵ Guaspari and Solovay let LTh also output the negations of the list members after outputting the members themselves, but that isn't necessary. An analogous point holds for case (ii), where Th outputs a *negated* Rosser sentence.

⁶ I do not however have a proof for this claim. The same result as in [1] is given in Smoryński's [4, Lemma 3.7, p. 295] for the usual proof predicate, but he doesn't have a similar problem because his definition of "standard proof predicate" [4, p. 279] already includes provable equivalence to $\text{Bew}(\mathbf{f})$, and thus our properties $(+^*)$ and $(+')$ (see p. 66).

Definition 1. $\text{Th}(\mathbf{f})$ is a *standard proof predicate* (SPP) iff $\text{Th}(\mathbf{f})$ is a Σ_1^0 -numeration of the theorems of \mathbf{PA} in \mathbf{PA} — i.e., there is a Δ_1^0 -formula $\text{PF}(\mathbf{p}, \mathbf{f})$ such that $\text{Th}(\mathbf{f}) = \exists \mathbf{p} \text{PF}(\mathbf{p}, \mathbf{f})$ and for all $\mathcal{L}_{\mathbf{PA}}$ -sentences φ ,

$$(\text{SPP1}) \quad \mathbb{N} \models \text{Th}(\ulcorner \varphi \urcorner) \Leftrightarrow \mathbf{PA} \vdash \varphi$$

holds — satisfying, for all sentences φ and ψ :

$$(\text{SPP2}) \quad \mathbf{PA} \vdash \text{Th}(\ulcorner \varphi \rightarrow \psi \urcorner) \wedge \text{Th}(\ulcorner \varphi \urcorner) \rightarrow \text{Th}(\ulcorner \psi \urcorner),$$

and for every Σ_1^0 -sentence σ :

$$(\text{SPP3}) \quad \mathbf{PA} \vdash \sigma \rightarrow \text{Th}(\ulcorner \sigma \urcorner).$$

Every SPP $\text{Th}(\mathbf{f})$ fulfills the Bernays–Löb derivability conditions, including, for all sentences φ :

$$(\text{DC1}) \quad \mathbf{PA} \vdash \varphi \Rightarrow \mathbf{PA} \vdash \text{Th}(\ulcorner \varphi \urcorner).^7$$

The extra property (+) for SPP's $\text{Th}(\mathbf{f})$, as given in [1, p. 97], is the following:

$$(+)\quad \mathbf{PA} \text{ proves: “}\{\mathbf{f} \mid \text{Th}(\mathbf{f})\} \text{ is closed under tautological consequence and contains all true } \Sigma_1^0\text{-sentences”},$$

where “true” is formalized by using the ‘usual’ truth predicate for Σ_1^0 -formulas.

To make this more explicit, let $\text{Taut}(\mathbf{f})$ be a Δ_1^0 -formula expressing that \mathbf{f} is (the Gödel number of) a tautological $\mathcal{L}_{\mathbf{PA}}$ -sentence, and let $\Sigma\text{True}(\mathbf{f})$ be Σ_1^0 , saying: “ \mathbf{f} is a true Σ_1^0 -sentence”, where a sentence's being true is analyzed as its being satisfied by some variable assignment.⁸ Furthermore, let $\text{Cond}(\mathbf{f}, \mathbf{g}, \mathbf{c})$ be a Δ_1^0 -pterm w.r.t. \mathbf{c} which characterizes \mathbf{c} as the material conditional of the formulas \mathbf{f} and \mathbf{g} . For this \mathbf{c} we subsequently write, not “ $\text{cond}(\mathbf{f}, \mathbf{g})$ ”, but, more suggestively, “ $\mathbf{f} \ominus \mathbf{g}$ ”. Then (+) consists in

$$\begin{aligned} (+_1) \quad & \mathbf{PA} \vdash \forall \mathbf{f} [\text{Taut}(\mathbf{f}) \rightarrow \text{Th}(\mathbf{f})], \\ (+_2) \quad & \mathbf{PA} \vdash \forall \mathbf{f}, \mathbf{g} [\text{Th}(\mathbf{f} \ominus \mathbf{g}) \wedge \text{Th}(\mathbf{f}) \rightarrow \text{Th}(\mathbf{g})], \\ (+_3) \quad & \mathbf{PA} \vdash \forall \mathbf{f} [\Sigma\text{True}(\mathbf{f}) \rightarrow \text{Th}(\mathbf{f})]. \end{aligned}$$

Guaspari and Solovay's proof of Lemma 6.3, however, needs an SPP $\text{Th}(\mathbf{f})$ with a little more than (+). This having been brought to his attention, Bob Solovay suggested strengthening (+) so as to include provable closure of Th under *first-order* logical consequence. The modified property (+*) then has (+₂) and (+₃) as before and, in place of (+₁), the following:

$$(+'_1) \quad \mathbf{PA} \vdash \forall \mathbf{f} [\text{PLTh}(\mathbf{f}) \rightarrow \text{Th}(\mathbf{f})],$$

where $\text{PLTh}(\mathbf{f})$ is Σ_1^0 , saying that there is an $\mathcal{L}_{\mathbf{PA}}$ -proof for sentence \mathbf{f} in first-order predicate logic, that is, a proof in the language of Peano arithmetic not using any arithmetical axioms.⁹

As an anonymous referee has remarked, what we actually make use of in proving Lemma 6.3, besides (+₂) and (+₃), is merely a weak consequence of (+'_1). — To enhance readability, I introduce some further pterms and other notation: $\text{CTm}(\mathbf{t})$ is to be a Δ_1^0 -formula saying that \mathbf{t} is a constant term; “ $\text{num}(\mathbf{x})$ ” is used to denote (the Gödel number of) the numeral of the number \mathbf{x} , a constant term; “ $\text{subst}(\mathbf{f}, \mathbf{v}, \mathbf{t})$ ” refers to the result of substituting the term \mathbf{t} for the variable \mathbf{v} in the formula \mathbf{f} ; and “ $\text{subst}_2(\mathbf{f}; \mathbf{v}, \mathbf{v}'; \mathbf{t}, \mathbf{t}')$ ” is short for

$$\text{subst}(\text{subst}(\mathbf{f}, \mathbf{v}, \mathbf{t}), \mathbf{v}', \mathbf{t}').$$

Finally, “ $\mathbf{f} \ominus \mathbf{g}$ ” and “ $\neg \mathbf{f}$ ” are to stand for the biconditional of the formulas \mathbf{f} and \mathbf{g} and the negation of \mathbf{f} , respectively, defined in analogy to “ $\mathbf{f} \oplus \mathbf{g}$ ”. All of these pterms are Δ_1^0 . — Now, as condition (+') we require that over and above (+₁)–(+₃), the following holds:

⁷ This is because \mathbf{PA} is Σ -complete and thus $\mathbb{N} \models \text{Th}(\ulcorner \varphi \urcorner)$ implies $\mathbf{PA} \vdash \text{Th}(\ulcorner \varphi \urcorner)$; cf. [5, Section 8.2].

⁸ Cf. [6, Ch. 9] or [3, Sect. 23].

⁹ Conditions (+'_1), (+₂) and (DC1) together imply (+₃).

(+4) for all formulas φ and all variables x, x' ,

$$\mathbf{PA} \vdash \forall t, t' [\text{CTm}(t) \wedge \text{CTm}(t') \wedge \text{Th}(\ulcorner \forall x, x' \varphi \urcorner) \rightarrow \text{Th}(\text{subst}_2(\ulcorner \varphi \urcorner; \ulcorner x \urcorner, \ulcorner x' \urcorner; t, t'))].^{10}$$

We prove the central lemma (cf. [3, [27.8]]) using the property (+'). For good measure, we also sketch a proof based on (+1)–(+3) and (+4') (see fn. 10), which is more similar to the original proof in [1].

Lemma 1. *PA proves: “if the bell rings, $\{f \mid \text{Th}(f)\}$ is inconsistent, i. e., $\text{Th}(\ulcorner \perp \urcorner)$ holds”.*

Proof. Reason in **PA**. Suppose, e. g., that the bell rings in step m because $\text{thNum}(m)$ itself is already on the Rosser list. Let's call it “ r ” for “Rosser sentence”. That r is on the list must be because thNum has previously output (the Gödel number of) a sentence of the form

$$\rho \leftrightarrow \exists y [\text{LPF}(y, \ulcorner \neg \rho \urcorner) \wedge (\forall z \leq y) \neg \text{LPF}(z, \ulcorner \rho \urcorner)],$$

viz., a sentence saying: “ r is a Rosser sentence w.r.t. $\text{LTh}(f)$.” I call such a biconditional's right-hand side the “Rosser proposition for ρ ”. Its construction can be captured in a Δ_1^0 -formula $\text{RProp}(f, p)$, a pterm w.r.t. p , and I write “ $\text{rProp}(f)$ ” for such a p . So, thNum has output $r \oplus \text{rProp}(r)$. Thus we have $\text{Th}(r \oplus \text{rProp}(r))$ and, since r is $\text{thNum}(m)$, also $\text{Th}(r)$. These two facts yield $\text{Th}(\text{rProp}(r))$: by (+1) we have

$$\text{Th}([r \oplus \text{rProp}(r)] \oplus [r \oplus \text{rProp}(r)]),$$

and two applications of (+2) give us $\text{Th}(\text{rProp}(r))$. Abbreviating the cumbersome formula $\exists y [\text{LPF}(y, f) \wedge (\forall z \leq y) \neg \text{LPF}(z, g)]$ as “ $\text{Before}(f, g)$ ”,¹¹ and letting $n := \ominus r$, we can denote $\text{rProp}(r)$ somewhat more transparently as “ $\text{Before}(\dot{n}, \dot{r})$ ”.¹² So, we have

$$\text{Th}(\ulcorner \text{Before}(\dot{n}, \dot{r}) \urcorner), \tag{1}$$

which means that Th says: “there is an LTh -disproof of r which occurs before any LTh -proof of r .”

But it's just the other way round in fact: The bell rang because $\text{thNum}(m) (= r)$ is on the Rosser list, so, by construction, LTh at step m starts outputting the whole Rosser list, including r , excluding its negation n (which, by construction, can't be on the list). Assume that LTh hasn't yet output n by step m . (If LTh should have proved n before the bell rang, then it would have done so because thNum , and thus Th , endorsed n . Since $\text{Th}(r)$ holds as well, Th would then be inconsistent anyway.) So there is an LTh -proof of r which occurs before any LTh -disproof of r :

$$\exists y [\text{LPF}(y, r) \wedge (\forall z \leq y) \neg \text{LPF}(z, n)],$$

that is,

$$\text{Before}(r, n).$$

¹⁰ The corresponding dual requirement

(+4') for all formulas φ and all variables x, x' ,

$$\mathbf{PA} \vdash \forall t, t' [\text{CTm}(t) \wedge \text{CTm}(t') \wedge \text{Th}(\text{subst}_2(\ulcorner \varphi \urcorner; \ulcorner x \urcorner, \ulcorner x' \urcorner; t, t')) \rightarrow \text{Th}(\ulcorner \exists x, x' \varphi \urcorner)]$$

can do the job too. In both cases, allowing just one substitution is *not* enough. A more straightforward, but stronger, consequence of (+1') would be

$$\mathbf{PA} \vdash \forall f, v, t [\text{Fml}(f) \wedge \text{Var}(v) \wedge \text{CTm}(t) \wedge \text{Th}(\bigotimes v f) \rightarrow \text{Th}(\text{subst}(f, v, t))],$$

where $\text{Fml}(f)$ and $\text{Var}(v)$ are Δ_1^0 -formulas saying that, respectively, f is a formula and v is a variable, and “ $\bigotimes v f$ ” stands for the v -universalization of f . Again, the dual sentence works as well.

¹¹ This corresponds to the expression “ $\text{Th}_f(f) < \text{Th}_f(g)$ ” in [1].

¹² The dots above the variables indicate that this is shorthand for

$$\text{subst}_2(\ulcorner \text{Before}(n, r) \urcorner; \ulcorner \dot{n} \urcorner, \ulcorner \dot{r} \urcorner; \text{num}(n), \text{num}(r)),$$

in which r and n are still free.

To expose Th as inconsistent, we have to show that Th proves this too. In order to do so, we use Th 's provable Σ -completeness, $(+_3)$. The expression

$$\text{subst}_2(\ulcorner \text{Before}(x, n) \urcorner; \ulcorner x \urcorner, \ulcorner n \urcorner; \text{num}(x), \text{num}(n)) = \ulcorner \text{Before}(\dot{x}, \dot{n}) \urcorner$$

is (the Gödel number of) a true Σ_1^0 -sentence; hence we have

$$\Sigma\text{True}(\ulcorner \text{Before}(\dot{x}, \dot{n}) \urcorner),$$

which by $(+_3)$ implies

$$\text{Th}(\ulcorner \text{Before}(\dot{x}, \dot{n}) \urcorner). \quad (2)$$

However, $\ulcorner \text{Before}(\dot{n}, \dot{x}) \urcorner$ and $\ulcorner \text{Before}(\dot{x}, \dot{n}) \urcorner$ do not yet contradict each other obviously. To bring the inconsistency out into the open, we finally put to use $(+_4)$.

Let's step out of \mathbf{PA} for a moment. For simple arithmetical reasons \mathbf{PA} proves

$$\forall f, g [\text{Before}(f, g) \rightarrow \neg \text{Before}(g, f)],$$

and thus by (DC1) also

$$\text{Th}(\ulcorner \forall f, g [\text{Before}(f, g) \rightarrow \neg \text{Before}(g, f)] \urcorner). \quad (3)$$

Now jump back inside \mathbf{PA} . To obtain $\text{Th}(\ulcorner \perp \urcorner)$ from (1) and (2), we have to apply the generalization in (3) to our particular 'sentences' x and n . Clause $(+_4)$, but not $(+_1)$ – $(+_3)$ by themselves, can give us what we need, viz., provability of the corresponding substitution instance:

$$\text{Th}(\ulcorner \text{Before}(\dot{x}, \dot{n}) \rightarrow \neg \text{Before}(\dot{n}, \dot{x}) \urcorner).$$

An application of $(+_2)$ yields $\text{Th}(\ulcorner \neg \text{Before}(\dot{n}, \dot{x}) \urcorner)$, and the rest is easy: by $(+_1)$, we have

$$\text{Th}(\ulcorner \text{Before}(\dot{n}, \dot{x}) \rightarrow [\neg \text{Before}(\dot{n}, \dot{x}) \rightarrow \perp] \urcorner),$$

and two more applications of $(+_2)$ finally deliver $\text{Th}(\ulcorner \perp \urcorner)$.

(The proof using $(+_4')$ runs as follows: From (1) and (2) we get (provability of) their tautological consequence,

$$\text{Th}(\ulcorner \text{Before}(\dot{n}, \dot{x}) \wedge \text{Before}(\dot{x}, \dot{n}) \urcorner),$$

and $(+_4')$ yields

$$\text{Th}(\ulcorner \exists f, g [\text{Before}(f, g) \wedge \text{Before}(g, f)] \urcorner).$$

Analogously to (3), however, we have

$$\text{Th}(\ulcorner \neg \exists f, g [\text{Before}(f, g) \wedge \text{Before}(g, f)] \urcorner),$$

and thus we get $\text{Th}(\ulcorner \perp \urcorner)$ again.)

The other case, where the bell rings because thNum outputs the negation n of some Rosser sentence x on the list, is similar but easier. We have $\text{Th}(\ominus x)$ and $\text{Th}(x \oplus \text{rProp}(x))$. By closure w.r.t. tautological consequence we get $\text{Th}(\ominus \text{rProp}(x))$, i.e.,

$$\text{Th}(\ulcorner \neg \text{Before}(\dot{n}, \dot{x}) \urcorner).$$

But in fact $\text{Before}(n, x)$ is true, because LTh outputs $n = \ominus x$ in or soon after step m , when x hasn't yet been output (or else Th is inconsistent anyway). As a Σ_1^0 -sentence, $\text{Before}(n, x)$ is thus provable:

$$\text{Th}(\ulcorner \text{Before}(\dot{n}, \dot{x}) \urcorner).$$

Here, the contradiction in Th 's theorems is obvious and so we don't need $(+_4)$ — nor $(+_4')$. As before, $(+_1)$ yields

$$\text{Th}(\ulcorner \text{Before}(\dot{n}, \dot{x}) \rightarrow [\neg \text{Before}(\dot{n}, \dot{x}) \rightarrow \perp] \urcorner),$$

and by twice applying $(+_2)$ we once more get $\text{Th}(\ulcorner \perp \urcorner)$. \square

Acknowledgements

I am grateful to Ulf Friedrichsdorf and Bob Solovay for their generous help, and to an anonymous referee for several useful suggestions.

References

- [1] D. Guaspari, R.M. Solovay, Rosser sentences, *Annals of Mathematical Logic* 16 (1) (1979) 81–99.
- [2] G. Boolos, *The Logic of Provability*, Cambridge University Press, Cambridge, New York, Oakleigh, 1993.
- [3] C. von Bülow, Beweisbarkeitslogik: Gödel, Rosser, Solovay, in: *Logische Philosophie*, vol. 15, Logos, Berlin, 2006.
- [4] C. Smoryński, Self-Reference and Modal Logic, in: *Universitext*, Springer, New York, Berlin, Heidelberg, Tokyo, 1985.
- [5] J.R. Shoenfield, *Mathematical Logic*, in: *Addison-Wesley Series in Logic*, Addison-Wesley, Reading, MA, Menlo Park, CA, London, Don Mills, Ontario, 1967.
- [6] R. Kaye, *Models of Peano Arithmetic*, Oxford University Press, Oxford, 1991.